

CSE574
Introduction
to Machine
Learning

Jue Guo

CSE574 Introduction to Machine Learning

Adversarial Attack: An Overview

Jue Guo

University at Buffalo

April 1, 2024

What are
adversarial
attacks?

The surprising
findings by
Szegedy (2013)
and Goodfellow
(2014)

Example of
attacks

Physical Attacks

Basic
Terminologies

Defining attacks

Multi-class
Problem

Three forms of
attack

Objective
function and
constraint sets

Outline

CSE574
Introduction
to Machine
Learning

Jue Guo

What are
adversarial
attacks?

The surprising
findings by
Szegedy (2013)
and Goodfellow
(2014)

Example of
attacks

Physical Attacks

Basic
Terminologies

Defining attacks

Multi-class
Problem

Three forms of
attack

Objective
function and
constraint sets

1 What are adversarial attacks?

- The surprising findings by Szegedy (2013) and Goodfellow (2014)
- Example of attacks
- Physical Attacks

2 Basic Terminologies

- Defining attacks
- Multi-class Problem
- Three forms of attack
- Objective function and constraint sets

Why?

CSE574
Introduction
to Machine
Learning

Jue Guo

What are
adversarial
attacks?

The surprising
findings by
szegedy (2013)
and goodfellow
(2014)

Example of
attacks

Physical Attacks

Basic
Terminologies

Defining attacks

Multi-class
Problem

Three forms of
attack

Objective
function and
constraint sets

- Robustness = easiness to fail when input is perturbed. Perturbation can be in any kind. Robustness machine learning is a very rich topic.
- We will look at something very narrow, called **adversarial robustness**, also known as robustness against **attacks**.
- Adversarial attack is a very **hot** topic, as of today. We should not over-emphasize its importance. There are many other important problems.

Adversarial Attack Example: FGSM

CSE574
Introduction
to Machine
Learning

Jue Guo

What are
adversarial
attacks?

The surprising
findings by
szegedy (2013)
and goodfellow
(2014)

Example of
attacks

Physical Attacks

Basic
Terminologies

Defining attacks

Multi-class
Problem

Three forms of
attack

Objective
function and
constraint sets

- It is not difficult to fool a classifier
- The perturbation could be perceptually not noticeable

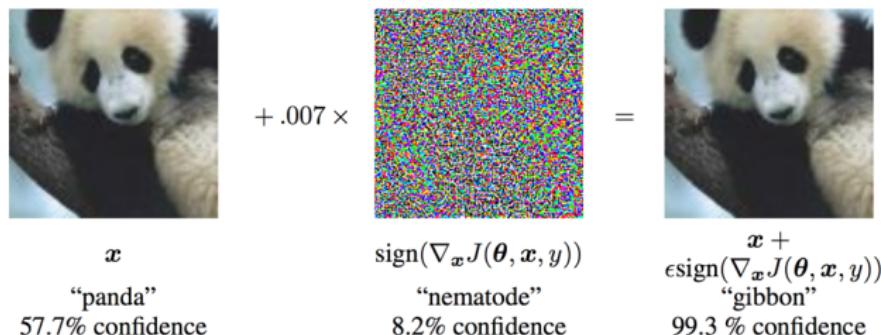


Figure: Goodfellow et al. "Explaining and Harnessing Adversarial Examples",
<https://arxiv.org/pdf/1412.6572.pdf>

Adversarial Attack Example: Szegedy's 2013 Paper

CSE574
Introduction
to Machine
Learning

Jue Guo

What are
adversarial
attacks?

The surprising
findings by
Szegedy (2013)
and Goodfellow
(2014)

Example of
attacks

Physical Attacks

Basic
Terminologies

Defining attacks

Multi-class
Problem

Three forms of
attack

Objective
function and
constraint sets

- This paper actually appears one year before Goodfellow's 2014 paper.

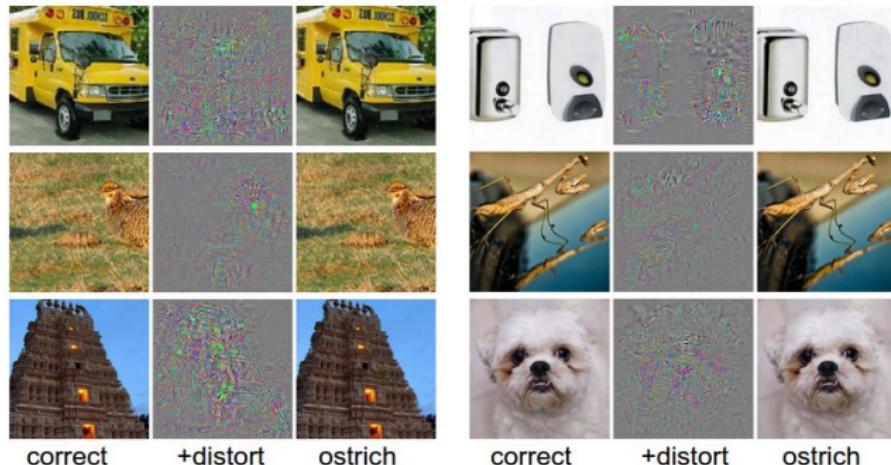


Figure: Szegedy et al. Intriguing properties of neural networks <https://arxiv.org/abs/1312.6199>

Adversarial Attack: Targeted Attack

CSE574
Introduction
to Machine
Learning
Jue Guo

What are
adversarial
attacks?

The surprising
findings by
Szegedy (2013)
and Goodfellow
(2014)

Example of
attacks

Physical Attacks

Basic
Terminologies

Defining attacks

Multi-class
Problem

Three forms of
attack

Objective
function and
constraint sets

■ Targeted Attack

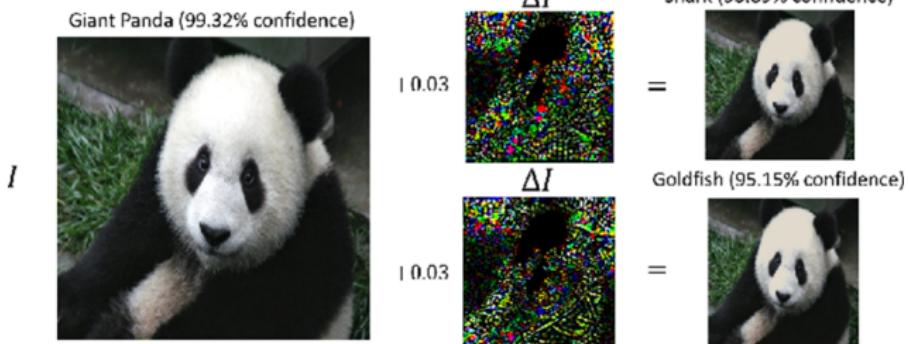


Figure: Adversarial Examples Detection in Deep Networks with Convolutional Filter Statistics,
<https://arxiv.org/abs/1612.07767>

Adversarial Attack Example: One Pixel

CSE574
Introduction
to Machine
Learning

Jue Guo

What are
adversarial
attacks?

The surprising
findings by
Szegedy (2013)
and Goodfellow
(2014)

Example of
attacks

Physical Attacks

Basic
Terminologies

Defining attacks

Multi-class
Problem

Three forms of
attack

Objective
function and
constraint sets

■ One-pixel Attack



SHIP
CAR(99.7%)



HORSE
FROG(99.9%)



DEER
AIRPLANE(85.3%)



DEER
DOG(86.4%)



HORSE
DOG(70.7%)



DOG
CAT(75.5%)



BIRD
FROG(86.5%)



BIRD
FROG(88.8%)

Figure: One pixel attack for fooling deep neural networks <https://arxiv.org/abs/1710.08864>

Adversarial Attack Example: Patch

CSE574
Introduction
to Machine
Learning

Jue Guo

What are
adversarial
attacks?

The surprising
findings by
Szegedy (2013)
and Goodfellow
(2014)

Example of
attacks

Physical Attacks

Basic
Terminologies
Defining attacks

Multi-class
Problem

Three forms of
attack

Objective
function and
constraint sets

■ Adding a patch



African-Elephant (92.8%) → Baseball (90.7%)



Sports Car (92.8%) → Shih-Tzu (90.7%)



Brown Bear (87.9%) → Tree Frog (82.7%)



Minivan (90.7%) → Tree Frog (86.4%)

Figure: LAVAN: Localized and visible Adversarial Noise, <https://arxiv.org/abs/1801.02608>

Adversarial Attack Example: Stop Sign

CSE574
Introduction
to Machine
Learning

Jue Guo

What are
adversarial
attacks?

The surprising
findings by
Szegedy (2013)
and Goodfellow
(2014)

Example of
attacks

Physical Attacks

Basic
Terminologies

Defining attacks

Multi-class
Problem

Three forms of
attack

Objective
function and
constraint sets

■ The Michigan / Berkely Stop Sign

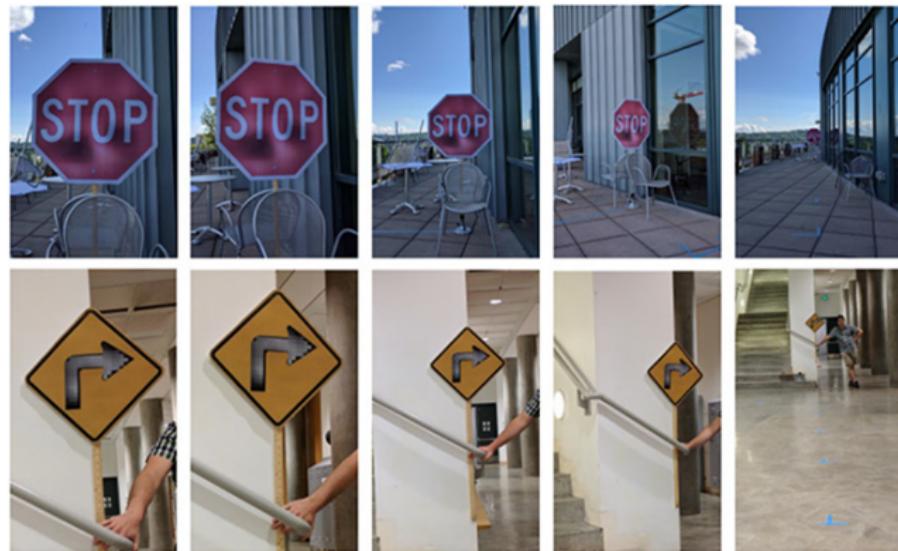


Figure: Robust Physical-World Attacks on Deep Learning Models

<https://arxiv.org/abs/1707.08945>

Adversarial Attack Example: Turtle

CSE574
Introduction
to Machine
Learning

Jue Guo

What are
adversarial
attacks?

The surprising
findings by
Szegedy (2013)
and Goodfellow
(2014)

Example of
attacks

Physical Attacks

Basic
Terminologies

Defining attacks

Multi-class
Problem

Three forms of
attack

Objective
function and
constraint sets

■ The MIT 3D Turtle



■ classified as turtle ■ classified as rifle ■ classified as other

Figure: Synthesizing Robust Adversarial Examples <https://arxiv.org/pdf/1707.07397.pdf>
<https://www.youtube.com/watch?v=YXy6oxiInoA>

Adversarial Attack Example: Glass

CSE574
Introduction
to Machine
Learning

Jue Guo

What are
adversarial
attacks?

The surprising
findings by
Szegedy (2013)
and Goodfellow
(2014)

Example of
attacks

Physical Attacks

Basic
Terminologies

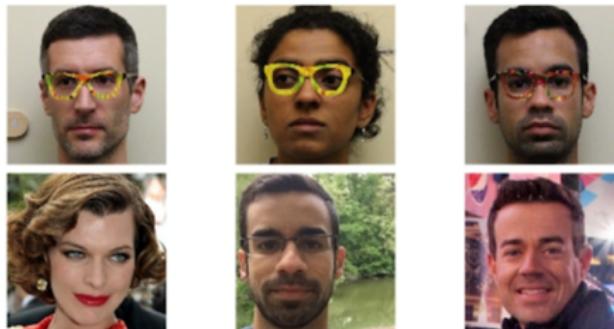
Defining attacks

Multi-class
Problem

Three forms of
attack

Objective
function and
constraint sets

■ CMU Glass



Input

Recognized Person

Sharif, M., Bhagavatula, S., Bauer, L., & Reiter, M. K. (2016, October). Accessorize to a crime: Real and stealthy attacks on state-of-the-art face recognition. In *Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security* (pp. 1528-1540). ACM.

Figure: Accessorize to a Crime: Real and Stealthy Attacks on State-of-the-Art Face Recognition
<https://www.cs.cmu.edu/~sbhagava/papers/face-rec-ccs16.pdf> <https://www.archive.ece.cmu.edu/~lbauer/proj/advm1.php>

Definition: Additive Adversarial Attack

CSE574
Introduction
to Machine
Learning

Jue Guo

What are
adversarial
attacks?

The surprising
findings by
Szegedy (2013)
and Goodfellow
(2014)

Example of
attacks

Physical Attacks

Basic
Terminologies

Defining attacks

Multi-class
Problem

Three forms of
attack

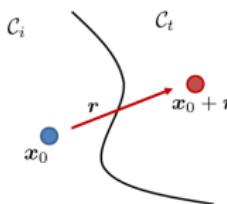
Objective
function and
constraint sets

Additive Adversarial Attack

Let $x_0 \in \mathbb{R}^d$ be a data point belong to class \mathcal{C}_i . Define a target class \mathcal{C}_t . An additive adversarial attack is an addition of a perturbation $r \in \mathbb{R}^d$ such that the perturbed data

$$x = x_0 + r$$

is misclassified as \mathcal{C}_t .



Definition: General Adversarial Attack

CSE574
Introduction
to Machine
Learning

Jue Guo

What are
adversarial
attacks?

The surprising
findings by
Szegedy (2013)
and Goodfellow
(2014)

Example of
attacks

Physical Attacks

Basic
Terminologies

Defining attacks

Multi-class
Problem

Three forms of
attack

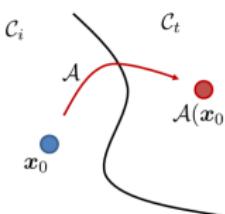
Objective
function and
constraint sets

General Adversarial Attack

Let $x_0 \in \mathbb{R}^d$ be a data point belonging to class \mathcal{C}_i . Define a target class \mathcal{C}_t . An adversarial attack is a mapping $\mathcal{A} : \mathbb{R}^d \rightarrow \mathbb{R}^d$ such that the perturbed data

$$x = \mathcal{A}(x_0)$$

is misclassified as \mathcal{C}_t .



Multi-class Problem

CSE574
Introduction
to Machine
Learning
Jue Guo

What are
adversarial
attacks?

The surprising
findings by
Szegedy (2013)
and Goodfellow
(2014)

Example of
attacks

Physical Attacks

Basic
Terminologies

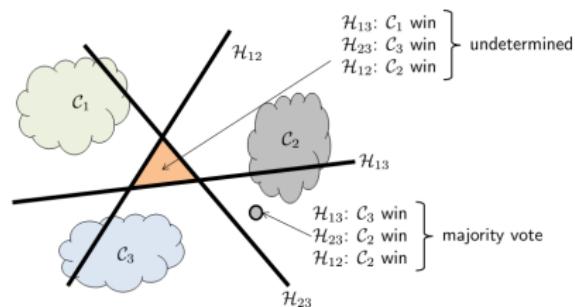
Defining attacks

Multi-class
Problem

Three forms of
attack

Objective
function and
constraint sets

Approach 1: One-on-One



- Class i vs. Class j
- Give me a point, check which class has more votes
- There is an undetermined region

The Multi-Class Problem

CSE574
Introduction
to Machine
Learning

Jue Guo

What are
adversarial
attacks?

The surprising
findings by
Szegedy (2013)
and Goodfellow
(2014)

Example of
attacks

Physical Attacks

Basic
Terminologies

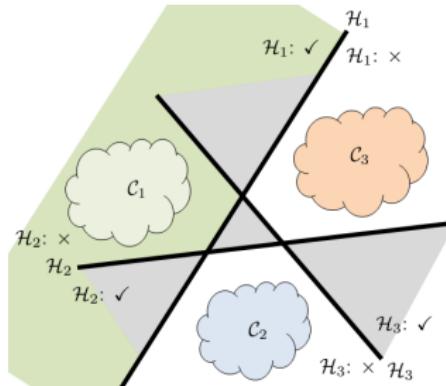
Defining attacks

Multi-class
Problem

Three forms of
attack

Objective
function and
constraint sets

Approach 2: One-on-All



- Class i not Class i
- Give me a point, check which class has no conflict
- There are undetermined regions

The Multi-Class Problem

CSE574
Introduction
to Machine
Learning

Jue Guo

What are
adversarial
attacks?

The surprising
findings by
Szegedy (2013)
and Goodfellow
(2014)

Example of
attacks

Physical Attacks

Basic
Terminologies

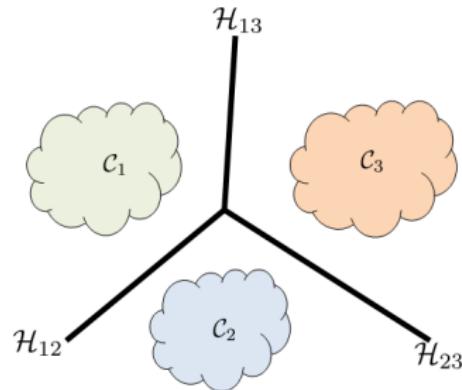
Defining attacks

Multi-class
Problem

Three forms of
attack

Objective
function and
constraint sets

Approach 3: Linear Machine



- Every point in the space gets assigned a class.
- You give me x , I compute $g_1(x), g_2(x), \dots, g_K(x)$
- If $g_i(x) \geq g_j(x)$ for all $j \neq i$, then x belongs to class i

Correct Classification

CSE574
Introduction
to Machine
Learning

Jue Guo

What are
adversarial
attacks?

The surprising
findings by
Szegedy (2013)
and Goodfellow
(2014)

Example of
attacks

Physical Attacks

Basic
Terminologies

Defining attacks

Multi-class
Problem

Three forms of
attack

Objective
function and
constraint sets

- We are mostly interested in the linear machine problem.
- Let us try to simplify the notation. The statement:

If $g_i(x) \geq g_j(x)$ for all $j \neq i$, then x belongs to class i

is equivalent to (asking everyone to be less than 0)

$$g_1(x) - g_i(x) \leq 0$$

⋮

$$g_k(x) - g_i(x) \leq 0$$

and is also equivalent to (asking the worst guy to be less than 0)

$$\max_{j \neq i} \{g_j(x)\} - g_i(x) \leq 0$$

- Therefore, if I want to launch an **adversarial attack**, I want to move you to class t :

$$\max_{j \neq t} \{g_j(x)\} - g_t(x) \leq 0$$

Our Approach

CSE574
Introduction
to Machine
Learning

Jue Guo

What are
adversarial
attacks?

The surprising
findings by
Szegedy (2013)
and Goodfellow
(2014)

Example of
attacks

Physical Attacks

Basic
Terminologies

Defining attacks

Multi-class
Problem

Three forms of
attack

Objective
function and
constraint sets

Here is what we are going to do

- First, we will preview the three **equivalent forms** of attack:
 - Minimum Distance Attack: Minimize the perturbation magnitude while accomplishing the attack objective.
 - Maximum Loss Attack: Maximize the training loss while ensuring perturbation is controlled.
 - Regularization-based Attack: Use regularization to control the amount of perturbation.
- Then, we will try to understand the **geometry** of the attacks.
- We will look at the **linear classifier** case to gain insights.

Minimum Distance Attack

CSE574
Introduction
to Machine
Learning

Jue Guo

What are
adversarial
attacks?

The surprising
findings by
Szegedy (2013)
and Goodfellow
(2014)

Example of
attacks

Physical Attacks

Basic
Terminologies

Defining attacks

Multi-class
Problem

Three forms of
attack

Objective
function and
constraint sets

Minimum Distance Attack

The **minimum distance attack** finds a perturbed data x by solving the optimization

$$\begin{aligned} & \underset{x}{\text{minimize}} && \|x - x_0\| \\ & \text{subject to} && \max_{j \neq t} \{g_j(x)\} - g_t(x) \leq 0 \end{aligned}$$

where $\|\cdot\|$ can be any norm specified by the user.

- I want to make you to class C_t
- So the constraint needs to be satisfied.
- But I also want to minimize the attack strength. This gives the objective.

Maximum Loss Attack

Maximum Loss Attack

The **maximum loss attack** finds a perturbed data x by solving the optimization

$$\begin{aligned} & \underset{x}{\text{maximize}} && g_t(x) - \max_{j \neq t} \{g_j(x)\} \\ & \text{subject to} && \|x - x_0\| \leq \eta \end{aligned}$$

where $\|\cdot\|$ can be any norm specified by the user, and $\eta > 0$ denotes the attack strength.

- I want to bound my attack $\|x - x_0\| \leq \eta$
- I want to make $g_t(x)$ as big as possible
- So I want to maximize $g_t(x) - \max_{j \neq t} \{g_j(x)\}$
- This is equivalent to

$$\begin{aligned} & \underset{x}{\text{minimize}} && \max_{j \neq t} \{g_j(x)\} - g_t(x) \\ & \text{subject to} && \|x - x_0\| \leq \eta \end{aligned}$$

Regularization-based Attack

CSE574
Introduction
to Machine
Learning

Jue Guo

What are
adversarial
attacks?

The surprising
findings by
Szegedy (2013)
and Goodfellow
(2014)

Example of
attacks

Physical Attacks

Basic
Terminologies

Defining attacks

Multi-class
Problem

Three forms of
attack

Objective
function and
constraint sets

Regularization-based Attack

The regularization-based attack finds a perturbed data x by solving the optimization

$$\underset{x}{\text{minimize}} \|x - x_0\| + \lambda \left(\max_{j \neq t} \{g_j(x)\} - g_t(x) \right)$$

where $\|\cdot\|$ can be any norm specified by the user, and $\lambda > 0$ is a regularization parameter.

- Combine the two parts via regularization
- By adjusting $(\epsilon, \eta, \lambda)$, all three will give the same optimal value.

Understanding the Geometry: Objective Function

CSE574
Introduction
to Machine
Learning

Jue Guo

What are
adversarial
attacks?

The surprising
findings by
Szegedy (2013)
and Goodfellow
(2014)

Example of
attacks

Physical Attacks

Basic

Terminologies

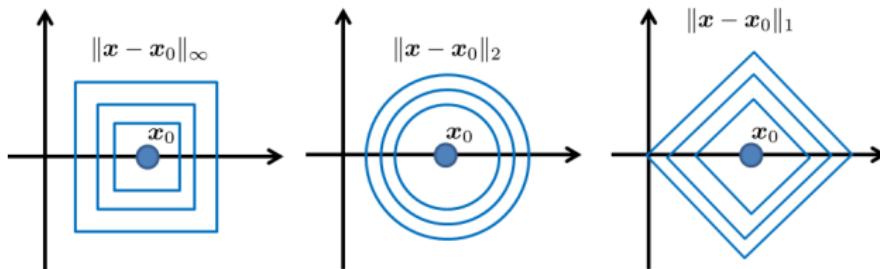
Defining attacks

Multi-class

Problem

Three forms of
attack

Objective
function and
constraint sets



- ℓ_0 -norm: $\varphi(x) = \|x - x_0\|_0$, which gives the most sparse solution. Useful when we want to limit the number of attack pixels.
- ℓ_1 -norm: $\varphi(x) = \|x - x_0\|_1$, which is a convex surrogate of the ℓ_0 .
- ℓ_∞ -norm: $\varphi(x) = \|x - x_0\|_\infty$, which minimizes the maximum element of the perturbation.

Understanding the Geometry: Constraint

CSE574
Introduction
to Machine
Learning

Jue Guo

What are
adversarial
attacks?

The surprising
findings by
Szegedy (2013)
and Goodfellow
(2014)

Example of
attacks

Physical Attacks

Basic
Terminologies
Defining attacks

Multi-class
Problem

Three forms of
attack

Objective
function and
constraint sets

- The constraint set is

$$\Omega = \left\{ x \mid \max_{j \neq t} \{g_j(x)\} - g_t(x) \leq 0 \right\}$$

- We can write Ω as

$$\Omega = \left\{ x \begin{array}{l} | g_1(x) - g_t(x) \leq 0 \\ | g_2(x) - g_t(x) \leq 0 \\ | \vdots \\ | g_k(x) - g_t(x) \leq 0 \end{array} \right\}$$

- Remark: If you want to replace \max by i^* , then i^* is a function of x :

$$\Omega = \left\{ x \mid g_{i^*(x)}(x) - g_t(x) \leq 0 \right\}$$

Understanding the Geometry: Constraint

CSE574
Introduction
to Machine
Learning

Jue Guo

What are
adversarial
attacks?

The surprising
findings by
Szegedy (2013)
and Goodfellow
(2014)

Example of
attacks

Physical Attacks

Basic
Terminologies

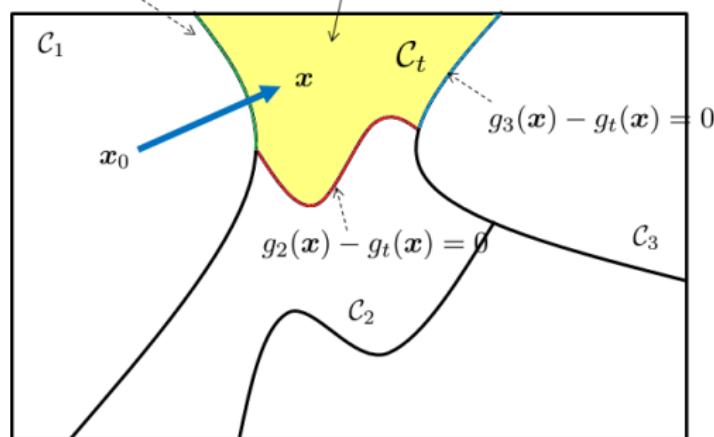
Defining attacks

Multi-class
Problem

Three forms of
attack

Objective
function and
constraint sets

$$g_1(\mathbf{x}) - g_t(\mathbf{x}) = 0 \quad \Omega = \left\{ \mathbf{x} \mid \max_{j \neq t} \{g_j(\mathbf{x})\} - g_t(\mathbf{x}) \leq 0 \right\}$$



Linear Classifier

CSE574
Introduction
to Machine
Learning

Jue Guo

What are
adversarial
attacks?

The surprising
findings by
Szegedy (2013)
and Goodfellow
(2014)

Example of
attacks

Physical Attacks

Basic
Terminologies

Defining attacks

Multi-class
Problem

Three forms of
attack

Objective
function and
constraint sets

- Let us take a closer look at the linear case.
 - Each discriminant function takes the form

$$g_i(x) = w_i^T x + w_{i,0}$$

- The decision boundary between the i -th class and the t -th class is therefore
 - The constraint set Ω is

$$g(x) = (w_i - w_t)^T x + w_{i,0} - w_{t,0} = 0$$

$$\begin{bmatrix} w_1^T - w_t^T \\ \vdots \\ w_{t-1}^T - w_t^T \\ w_{t+1}^T - w_t^T \\ \vdots \\ w_k^T - w_t^T \end{bmatrix} x + \begin{bmatrix} w_{1,0} - w_{t,0} \\ \vdots \\ w_{t-1,0} - w_{t,0} \\ w_{t+1,0} - w_{t,0} \\ \vdots \\ w_{k,0} - w_{t,0} \end{bmatrix} \leq \mathbf{0} \Leftrightarrow A^T x \leq b$$

Linear Classifier

CSE574
Introduction
to Machine
Learning

Jue Guo

What are
adversarial
attacks?

The surprising
findings by
Szegedy (2013)
and Goodfellow
(2014)

Example of
attacks

Physical Attacks

Basic
Terminologies

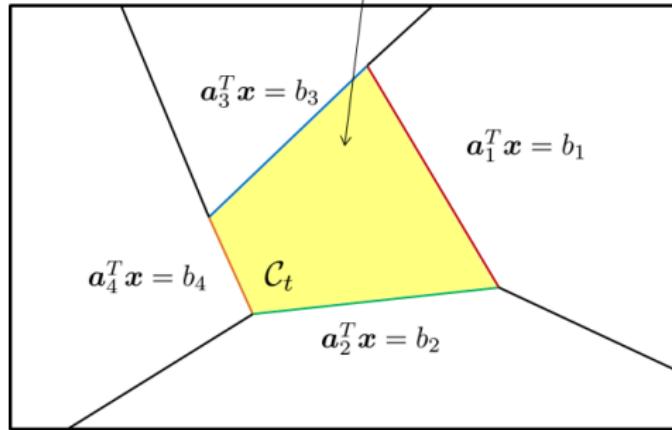
Defining attacks

Multi-class
Problem

Three forms of
attack

Objective
function and
constraint sets

$$\Omega = \{x \mid A^T x \leq b\}$$



- You can show $\Omega = \{A^T x \leq b\}$ is convex.
- But the complement $\Omega^c = \{A^T x > b\}$ is not convex.
 - So targeted attack is easier to analyze than untargeted attack.

Attack: The Simplest Example

CSE574
Introduction
to Machine
Learning

Jue Guo

What are
adversarial
attacks?

The surprising
findings by
Szegedy (2013)
and Goodfellow
(2014)

Example of
attacks

Physical Attacks

Basic
Terminologies
Defining attacks
Multi-class
Problem

Three forms of
attack

Objective
function and
constraint sets

The optimization is:

$$\begin{aligned} & \underset{x}{\text{minimize}} && \|x - x_0\| \\ & \text{subject to} && \max_{j \neq t} \{g_j(x)\} - g_t(x) \leq 0 \end{aligned}$$

- Suppose we use ℓ_2 -norm, and consider linear classifiers, then the attack is given by

$$\underset{x}{\text{minimize}} \|x - x_0\|^2 \text{ subject to } A^T x \leq b$$

- This is a quadratic programming problem.
- We will discuss how to solve this problem analytically.

Summary

CSE574
Introduction
to Machine
Learning

Jue Guo

What are
adversarial
attacks?

The surprising
findings by
Szegedy (2013)
and Goodfellow
(2014)

Example of
attacks

Physical Attacks

Basic
Terminologies

Defining attacks

Multi-class
Problem

Three forms of
attack

Objective
function and
constraint sets

- Adversarial attack is a universal phenomenon for **any** classifier.
- Attacking deep networks are popular because people think that they are unbeatable.
- There is really nothing too magical behind adversarial attack.
 - All attacks are based on one of the forms of attacks.
 - Deep networks are trickier, as we will see, because the internal model information is not easy to extract.
 - We will learn the basic principles of attacks, and try to gain insights from linear models.